

Ask-Me-Why: Model-agnostic Explainability Framework

Jonathan Lai and Hegler Tissot (Advisor)

Drexel University



PROBLEM & MOTIVATION

The Challenge

Current explainability methods require access to model internals. What if we could explain ANY model's predictions using only the data?

The Black Box Challenge

Healthcare: Why was this patient diagnosed with condition X?

Finance: Why was this loan application rejected?

Law: Why was this case flagged as high-risk?

Current Limitations

Method	Limitation
SHAP	Requires model access; Different models → different explanations
LIME	Model-specific; Unstable for mixed data types
XGBoost	Tied to one model; Often global rather than instance-specific

OBJECTIVES

- Explain ANY model without access using neighborhood analysis - Identify instance patterns & provide complementary insights to SHAP/LIME/XGBoost
- Enable regulatory compliance for healthcare, finance, law
- Support mixed data types & scale efficiently
- Empower domain experts with interpretable explanations
- Provide model prediction explanations

Key Insight

Similar instances in embedding space should behave similarly. When they don't, those differences explain the prediction.

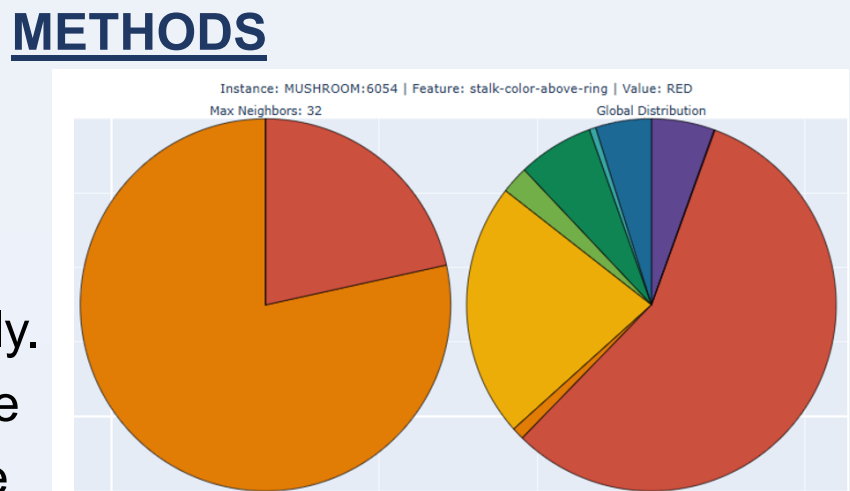


Figure 1: Local vs. Global

Core Innovation

Weighted Cosine Difference (WCD):

$$WCD(f, x) = \text{CosDiff}(P_{\text{local}}, P_{\text{global}}) \times \text{Uniqueness}(x_f)$$

Where:

P_local = Feature distribution in k-nearest neighbors

P_global = Feature distribution in entire dataset

$$\text{CosDiff} = 1 - \frac{P_{\text{local}} \cdot P_{\text{global}}}{\|P_{\text{local}}\| \cdot \|P_{\text{global}}\|}$$

$$\text{Uniqueness}_{\text{cat}}(x_f) = 1 - \frac{P_{\text{local}}(x_f)}{\sum P_{\text{local}}}$$

$$\text{Uniqueness}_{\text{cont}}(x_f) = \frac{|x_f - \mu_{\text{local}}|}{\sigma_{\text{local}}}$$

Where:

x_f = instance's feature value

μ_local = neighborhood mean

σ_local = neighborhood standard deviation

Framework:

Data → Embed → Analyze → Model → Explain

Algorithm Pipeline:

- Embed instances using unsupervised methods
- Build KD-Tree index for neighbor search
- Find k-nearest neighbors in embedding space
- Compare local vs global feature distributions
- Weight by instance uniqueness → WCD score
- Train model on same instances
- Provide model explanations

Technical Details

Feature Type	Missing Value Handling	Methods
Continuous	Treated as "Missing" bin; excluded from numeric bins	<ul style="list-style-type: none">Adaptive binning (0.25σ, 0.5σ, 1σ, 2σ, 4σ)Linear interpolation for bin assignmentDistance-weighted aggregation
Categorical	Treated as "Missing" category	<ul style="list-style-type: none">Frequency-based distributionsDistance-weighted counting

Datasets

Dataset	Instances	Features	Task	Characteristics
Mushroom	8,124	22 categorical	Binary	Colors & Structure
Genetic Disorder	22,383	6 cont, 36 cat	Multi-label	Age & lab patterns

Neighborhood Search:

KD-Tree indexing: $O(n \log n)$

$k \in \{16, 32, 64, 128, 256, 512, 1024\}$

Max radius: 2.0 (Euclidean)

Distance weighting:

$$w_i = \frac{r_{\text{max}} - d_i}{r_{\text{max}}}$$

d_i = distance to neighbor

r_max = maximum radius

Comparison Methods:

SHAP (TreeExplainer) - Shapley value approximation

LIME (Tabular) - Local linear surrogate

XGBoost - Global feature importance

Validation Approach:

2 diverse datasets with different characteristics

Spearman rank correlation for method comparison

Top-k feature agreement analysis

RESULTS

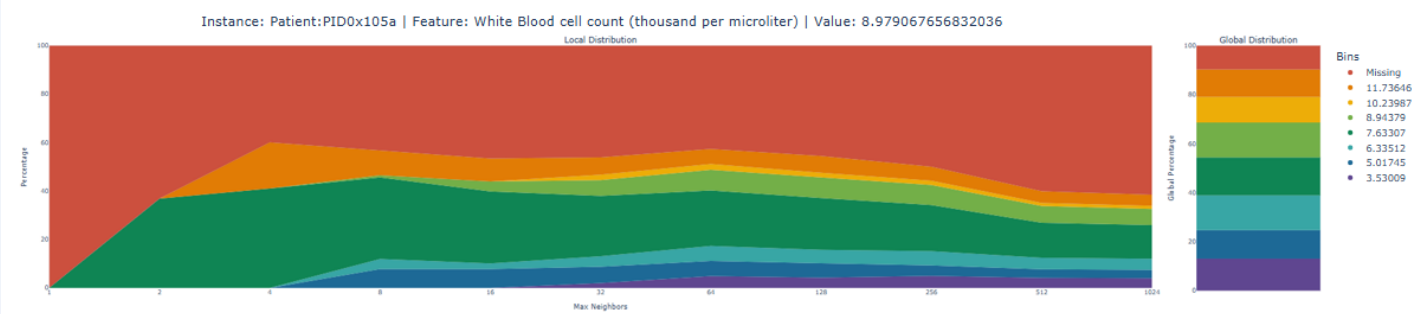


Figure 2. Local vs. Global White-Blood-Cell Distributions

feature	Rank	cos_diff	uniqueness	weighted_cos_diff_no_rmalized
White Blood cell count	1	0.41886	1.931932	0.33363

White Blood Cell count "neighborhood" shifts as we vary the number of nearest neighbors (k). The right-hand bar is the fixed global distribution.

Feature Ranking Comparison:

Table 1: Mushroom Dataset (128 neighbors)

Method Comparison	Spearman ρ	Interpretation
SHAP	0.016	Weak positive
LIME	0.450	Moderate positive
XGBoost	0.186	Weak positive

Table 2: Genetic Disorder (512 neighbors)

Method Comparison	Spearman ρ	Interpretation
SHAP	0.128	Weak positive
LIME	0.121	Weak positive
XGBoost	0.005	Neutral

Model Explanations

1. Mushroom Instance #3900: (#1 XGBoost: Odor)

Prediction: EDIBLE ✓ (99.97% confidence)

Top Features (ASK-Me-Why):

- Gill-color (BROWN)
- Stalk-color below ring (GRAY)
- Cap-color (BROWN)

2. Patient PID0x7005:

Prediction: Mitochondrial disorder ✓ (67.9%)

Top Features (ASK-Me-Why):

- Mother's age (33)
- Patient age (0)
- Father's age (52)

CONCLUSIONS

- Ask-Me-Why produces per-instance explanations **without** model access.
- Weighted Cosine Difference** highlights features both distributionally distinct and value-extreme.
- Complements SHAP, LIME & XGBoost by surfacing **different** but meaningful signals and comparing feature rankings.
- Scales to mixed data & large datasets via KD-Tree k-NN search.
- Provides visualization tools for further analysis

REFERENCES

- Ribeiro M.T., Singh S., Guestrin C. "Why Should I Trust You? Explaining the Predictions of Any Classifier." *KDD '16*, pp. 1135–1144.
- Lundberg S.M., Lee S.-I. "A Unified Approach to Interpreting Model Predictions." *NeurIPS '17*, pp. 4765–4774.
- Chen T., Guestrin C. "XGBoost: A Scalable Tree Boosting System." *KDD '16*, pp. 785–794.
- Linardatos P., Papastefanopoulos V., Kotsiantis S. "Explainable AI: A Review of Machine Learning Interpretability Methods." *Entropy* 23(1): 18, 2021.